

BAYESIAN CENTROID ESTIMATION FOR MOTIF DISCOVERY

BY LUIS CARVALHO*

Boston University

Biological sequences may contain patterns that are signal important biomolecular functions; a classical example is regulation of gene expression by transcription factors that bind to specific patterns in genomic promoter regions. In motif discovery we are given a set of sequences that share a common motif and aim to identify not only the motif composition, but also the binding sites in each sequence of the set. We present a Bayesian model that is an extended version of the model adopted by the Gibbs motif sampler, and propose a new centroid estimator that arises from a refined and meaningful loss function for binding site inference. We discuss the main advantages of centroid estimation for motif discovery, including computational convenience, and how its principled derivation offers further insights about the posterior distribution of binding site configurations. We also illustrate, using simulated and real datasets, that the centroid estimator can differ from the maximum a posteriori estimator.

1. Introduction. In motif discovery we are given a set of sequences that share a common motif and aim to identify the motif composition—the frequency of symbols for each position in the pattern—and the positions in each sequence where the motifs are. It is assumed that the motifs are significantly different, in composition, from sequence background. This problem has gained attention and relevance in the past 25 years mainly due to biological applications; a classical example is regulation of gene expression by transcription factors that bind to specific motifs in genomic promoter regions (MacIsaac and Fraenkel, 2006; GuhaThakurta, 2006; Sandve and Drablos, 2006). For this reason, we refer to the positions where the motifs are realized in the sequences as “binding sites”.

Due to its importance, hundreds of procedures have been proposed for motif discovery (Hu et al., 2005; Tompa et al., 2005). While some approaches seek to characterize motifs and their binding sites using dictionary methods that capture over-representation of words as evidence (Régner and Denise, 2004; Pavesi et al., 2004), it is common to represent motif compositions by a position weight matrix (Stormo, 2000) and specify a parametric model where sequences are generated conditionally on motif and background compositions and binding sites. Binding sites can then be regarded as missing data; parameters for the compositions can be estimated using expectation-maximization (Dempster et al., 1977) in a frequentist setup, as in MEME (Bailey and Elkan, 1995), or assigned a prior distribution in a Bayesian setup (Lawrence et al., 1993; Neuwald et al., 1995).

Following the Bayesian model from (Liu et al., 1995), we assume that there is only one motif of *fixed* length L and that sequences are generated conditionally independently according to a product multinomial model given binding site positions and motif and background compositions. Thus, for an alphabet \mathcal{S} , we define $\theta_0 = (\theta_{0,s})_{s \in \mathcal{S}}$ as background probabilities of generating each letter in \mathcal{S} and, for each position $i = 1, \dots, L$ in the motif, $\theta_i = (\theta_{i,s})_{s \in \mathcal{S}}$ as the probabilities of generating each letter at the i -th position in the motif. To simplify

*Supported by NSF grant DMS-1107067.

Keywords and phrases: Gibbs sampling, stochastic backtracking
AMS 2000 subject classifications: 62F15, 62P10, 65C05

the notation we denote $\Theta = (\theta_0, \theta_1, \dots, \theta_L)$. As in (Liu et al., 1995), we set a conjugate Dirichlet prior for Θ .

Product multinomial and product Dirichlet models are justified as a good working, first approximation based on position independence. There are many extensions to this model that consider DNA strand complementarity (Roth et al., 1998), a more informative Markov structure for the background composition (Liu et al., 2001), and an explicit representation of the number of binding sites per sequence (Thijs et al., 2002). However, since we will be discussing a new inferential procedure, we adopt an extended model that yields a feasible computational method while still retaining a realistic interpretation and allows us to focus the discussion on the proposed estimator.

Motif discovery is considered a hard problem since motifs are usually short relative to sequence length and have a composition that might be hard to distinguish from background (see, for instance, (Hu et al., 2005).) It is then imperative to rely on more refined, informative estimation methods that better glean information from the posterior distribution of binding site configurations. Discrete inferential methods with this goal have recently been proposed, including the median probability model of Barbieri and Berger (2004) and the centroid estimator (Ding et al., 2005; Carvalho and Lawrence, 2008). Centroid estimation, in particular, has been successfully used for motif discovery (Thompson et al., 2007), including models that account for sequence conservation (Newberg et al., 2007).

In this paper we present a Bayesian model for motif discovery on multiple sequences with multiple possible binding sites and formalize a new flavor of inference based on centroid estimation. As we will argue, the proposed estimator offers a good representative of the posterior space of binding site configurations; moreover, as a by-product of its derivation, we obtain informative summaries of the distribution of posterior mass. We start the discussion by addressing a simple case when there is only one sequence and we accept only one binding site; next we extend the presentation to include multiple binding sites; then, we treat the full case when Θ is random, in a fully Bayesian approach. Finally, we offer some concluding remarks and directions for future work in the last section.

2. One sequence, one binding site. Suppose we observe a sequence R , $|R| \doteq n$, and wish to infer the location of the only binding site Y , $Y \in \{1, \dots, n - L + 1\}$. Setting a non-informative prior on Y , $\mathbb{P}(Y) = (n - L + 1)^{-1}$, we have the posterior:

$$\mathbb{P}(Y | R, \Theta) = \frac{\mathbb{P}(R | Y, \Theta) \mathbb{P}(Y | \Theta)}{\sum_{\tilde{Y}=1}^{n-L+1} \mathbb{P}(R | \tilde{Y}, \Theta) \mathbb{P}(\tilde{Y} | \Theta)} = \frac{\mathbb{P}(R | Y, \Theta)}{\sum_{\tilde{Y}=1}^{n-L+1} \mathbb{P}(R | \tilde{Y}, \Theta)}.$$

The likelihood, as previously stated, follows a product multinomial distribution given Y :

$$\mathbb{P}(R | Y, \Theta) = \prod_{s \in \mathcal{S}} \prod_{j \in BG} \theta_{0,s}^{I(R_j=s)} \prod_{j=1}^L \theta_{j,s}^{I(R_{Y-j+1}=s)},$$

where $j \in BG$ means position j in background.

One traditional estimator is the MAP estimator,

$$\hat{Y}_M = \arg \max_{\tilde{Y}=1, \dots, n-L+1} \mathbb{P}(\tilde{Y} | R, \Theta),$$

but we argue for an estimator that accounts for differences in positions when comparing binding site configurations. Using Bayesian decision theory (Berger, 1985) we look for an

estimator that minimizes, on average, a more refined loss function H :

$$(1) \quad \hat{Y}_C = \arg \min_{\tilde{Y}=1, \dots, n-L+1} \mathbb{E}_{Y|R, \Theta} [H(\tilde{Y}, Y)].$$

We adopt a generalized Hamming loss H ,

$$H(\tilde{Y}, Y) = \sum_{i=1}^n h(l_i(\tilde{Y}), l_i(Y)),$$

where $l_i(Y)$ returns the “state” of position i : if i is a background position, $l_i(Y) = 0$, otherwise $l_i(Y) = Y - i + 1$, that is, $l_i(Y)$ returns the position in the motif. Loss function H compares configurations position-wise according to h , which in turn compares states. One option for h when Θ is known is a probability distance, the symmetric Kullback-Leibler distance,

$$h(i, j) = D_{KL}(\theta_i \parallel \theta_j) + D_{KL}(\theta_j \parallel \theta_i) = \sum_{s \in \mathcal{S}} \theta_{i,s} \log \frac{\theta_{i,s}}{\theta_{j,s}} + \theta_{j,s} \log \frac{\theta_{j,s}}{\theta_{i,s}},$$

for $i, j = 0, 1, \dots, L$.

It is, however, not common to have such an informed loss function. An alternative metric arises by simply allowing $\theta_{j,s} \doteq \theta_s \neq \theta_{0,s}$ for all $s \in \mathcal{S}$ and $j = 1, \dots, L$ in the motif. In this case, if $m(i) \doteq I(i > 0)$ indicates if state i is a motif state,

$$h(i, j) = h(m(i), m(j)) = I(m(i) \neq m(j)) \left[\sum_{s \in \mathcal{S}} \theta_s \log \frac{\theta_s}{\theta_{0,s}} + \theta_{0,s} \log \frac{\theta_{0,s}}{\theta_s} \right].$$

Since we are ultimately concerned with the argument of a minimum, as per Equation 1, we can define the loss function up to a shift and (positive) scale. Thus, for our inferential purposes it suffices to define $h(i, j) = I(m(i) \neq m(j))$ to obtain a loss H that accounts for overlap in binding sites. Such metric is commonly adopted to measure binding site level accuracy, as in the performance coefficients in (Pevzner et al., 2000; Hu et al., 2005; Tompa et al., 2005). From now on we will be focusing on this minimally informed loss function.

Estimator \hat{Y}_C is a *generalized centroid estimator*; for instance, if h is a common zero-one loss, $h(i, j) = I(i \neq j)$, H corresponds to Hamming loss, and thus \hat{Y}_C is the regular centroid estimator (Ding et al., 2005; Carvalho and Lawrence, 2008). As Carvalho and Lawrence (2008) argue, centroid estimators more effectively represent the space since they are closer to posterior means; in contrast, it can be shown that \hat{Y}_M arises from a zero-one loss function which yields the posterior mode (Besag, 1986).

Let us now derive more specific expressions for H and \hat{Y}_C . We first notice that if $|\tilde{Y} - Y| \geq L$ then the binding sites do not overlap and so $H(\tilde{Y}, Y) = 2 \sum_{j=1}^L h(j, 0) \doteq H^*$, the null overlap distance between two configurations. Alternatively, when $|\tilde{Y} - Y| < L$ then

$$(2) \quad H(\tilde{Y}, Y) = \sum_{j=1}^{|\tilde{Y}-Y|} h(j, 0) + \sum_{j=L-|\tilde{Y}-Y|+1}^L h(j, 0) + \sum_{j=1}^{L-|\tilde{Y}-Y|} h(j, j + |\tilde{Y} - Y|),$$

since the common backgrounds in \tilde{Y} and Y do not affect $H(\tilde{Y}, Y)$, the first two terms above account for the left and right “tails” where binding sites in one sequence are matched with

background in the other sequence, and the last term accounts for the overlap in binding sites. We also note that $H(\tilde{Y}, Y)$ is actually a function of $|\tilde{Y} - Y|$.

Instead of a loss function we can also define our estimator in terms of a *gain* function $G(\tilde{Y}, Y) \doteq 1 - H(\tilde{Y}, Y)/H^*$. Note that $0 \leq G(\tilde{Y}, Y) \leq 1$; in particular, when $|\tilde{Y} - Y| \geq L$ there is no gain, $G(\tilde{Y}, Y) = 0$, and if $\tilde{Y} = Y$ we have $G(\tilde{Y}, Y) = 1$. As a consequence, we can simply write $G(\tilde{Y}, Y) = I(|\tilde{Y} - Y| < L)(1 - H(\tilde{Y}, Y)/H^*)$ with H from Equation 2. Noting that G , like H , is also a function of $|\tilde{Y} - Y|$, we obtain the following characterization:

THEOREM 1. *The centroid estimator \hat{Y}_C is*

$$\hat{Y}_C = \arg \max_{\tilde{Y}=1, \dots, n-L+1} G(\tilde{Y}, \cdot) * \mathbb{P}(\cdot | R, \Theta),$$

a convolution between G and the posterior distribution on Y .

PROOF. The result follows directly from the definition in Equation 1:

$$\begin{aligned} \hat{Y}_C &= \arg \min_{\tilde{Y}=1, \dots, n-L+1} \mathbb{E}_{Y|R, \Theta} [H(\tilde{Y}, Y)] \\ &= \arg \max_{\tilde{Y}=1, \dots, n-L+1} \mathbb{E}_{Y|R, \Theta} [I(|\tilde{Y} - Y| < L)(1 - H(\tilde{Y}, Y)/H^*)] \\ &= \arg \max_{\tilde{Y}=1, \dots, n-L+1} \sum_{Y=\max\{1, \tilde{Y}-L+1\}}^{\min\{n-L+1, \tilde{Y}+L-1\}} G(\tilde{Y}, Y) \mathbb{P}(Y | R, \Theta) \\ &= \arg \max_{\tilde{Y}=1, \dots, n-L+1} G(\tilde{Y}, \cdot) * \mathbb{P}(\cdot | R, \Theta), \end{aligned}$$

as required. \square

When contrasted to \hat{Y}_M we can see the effect of having a higher resolution loss function: \hat{Y}_C gathers probability support from nearby, relative to H , binding site configurations instead of just picking the most likely configuration. The following example should give us some insight into this new estimator.

EXAMPLE 1. Consider the following sequence of length $n = 200$ from the nucleotide alphabet $\mathcal{S} = \{\text{A}, \text{C}, \text{G}, \text{T}\}$,

10	20	30	40	50
GCCACTTTCGGGCCCGTGTCTAACGCACCACGGGCTACGTGACGGTGTGG				
CTCTATACTGACGACGTGAACCAAGCTTTACTGAAGGACTTGCTGTTCCC				
CGACCCATTTCCTGCCAGAACCTCTGACCAGTGTCTAGGGCTATCGCCCCG				
TGATGTCTCATGGCGACGCGGAGGCGGTTGCTCGCCTCACTCCGTTCTG				

and a motif of length $L = 6$ with parameters Θ given by Table 1.

Figure 1 shows the conditional marginal posterior $\mathbb{P}(Y | R, \Theta)$ and the convolution $G * \mathbb{P}(\cdot | R, \Theta)$ used to obtain the centroid $\hat{Y}_C = 36$, binding at the subsequence TACGTG, close to the consensual motif. Note that since Θ is very informative the posterior profile has clear peaks and in this case $\hat{Y}_C = \hat{Y}_M$, the two estimators coincide.

TABLE 1

Background and motif compositions: background is assumed to be CG-rich, while the motif represents a canonical palindromic E-box, CACGTG (Murrea et al., 1989).

\mathcal{S}	θ_0	θ_1	θ_2	θ_3	θ_4	θ_5	θ_6
A	0.2	0.1	0.7	0.1	0.1	0.1	0.1
C	0.3	0.7	0.1	0.7	0.1	0.1	0.1
G	0.3	0.1	0.1	0.1	0.7	0.1	0.7
T	0.2	0.1	0.1	0.1	0.1	0.7	0.1

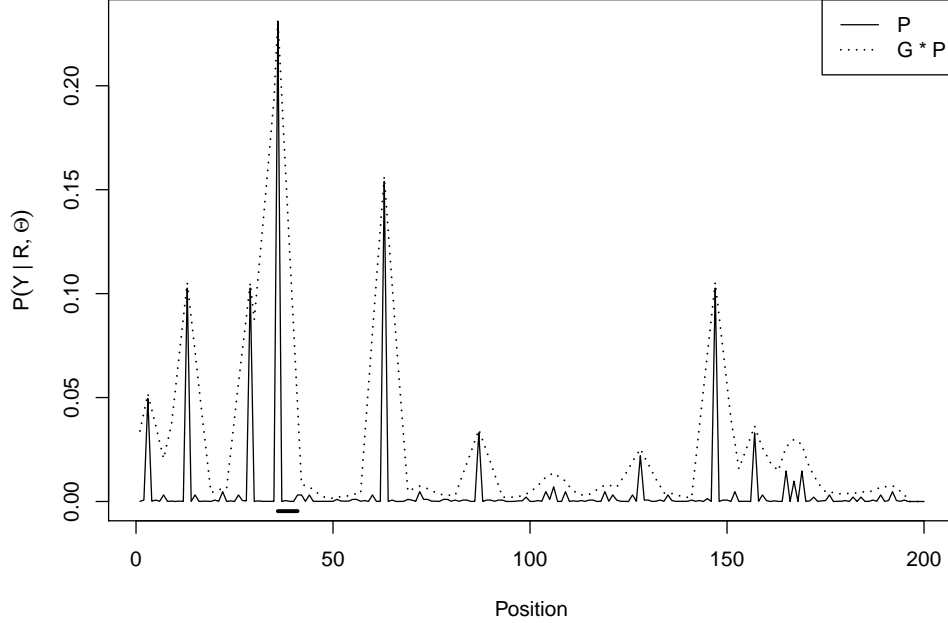


FIG 1. Conditional marginal probability distribution $\mathbb{P}(Y | R, \Theta)$ in solid line and convolution $G * \mathbb{P}(\cdot | R, \Theta)$ in dotted line. The black thick line close to the axis marks the binding site corresponding to the centroid \hat{Y}_C .

3. One sequence, multiple binding sites. We now allow for multiple binding sites by defining $Y = \{Y_k\}$ as the collection of binding sites Y_k . The likelihood is similar, but accounts for the multiple binding sites:

$$\mathbb{P}(R | Y, \Theta) = \prod_{s \in \mathcal{S}} \prod_{i \in BG} \theta_{0,s}^{I(R_i=s)} \prod_{k=1}^{|Y|} \prod_{i=1}^L \theta_{i,s}^{I(R_{Y_k+i-1}=s)}.$$

Given the “entropic” effect of possibly having many binding sites, we need to adopt a better prior for Y that takes into account the number of possible configurations for the binding sites. So, instead of naively electing $\mathbb{P}(Y) \propto 1$, we can explore a hierarchical structure: if $c(Y) = |Y|$, the number of binding sites in Y , we note that $\mathbb{P}(Y) = \mathbb{P}(Y, c(Y)) = \mathbb{P}(Y | c(Y))\mathbb{P}(c(Y))$ and set $\mathbb{P}(Y | c(Y)) \propto 1$ and $\mathbb{P}(c(Y)) \propto 1$ to obtain

$$\mathbb{P}(Y) = \mathbb{P}(Y | c(Y))\mathbb{P}(c(Y)) = \binom{n - c(Y)(L-1)}{c(Y)}^{-1} \cdot \frac{1}{C},$$

where $C \doteq \lfloor n/L \rfloor$ is the maximum number of binding sites in R .

Another, possibly more familiar, approach is to adopt a Markov chain with two states, background and motif, where the probability of transitioning to background, either from background or motif, and of starting at background is p . In this case we keep $\mathbb{P}(Y | c(Y))$ as before, but now

$$(3) \quad \mathbb{P}(c(Y)) \propto \binom{n - c(Y)(L - 1)}{c(Y)} p^{n - c(Y)L} (1 - p)^{c(Y)},$$

since there needs to be $c(Y)$ transitions to the motif state. This prior structure offers more flexibility through p : we can further set a hyperprior distribution on p , or specify it directly based on the expected number b of binding sites in the sequence; if n is large compared to b , as usual, then p should be close to one, $c(Y)$ is approximately Poisson with mean $n(1 - p)$ and thus $p \doteq 1 - b/n$ becomes a good candidate.

The posterior is then

$$\begin{aligned} \mathbb{P}(Y | R, \Theta) &= \frac{\mathbb{P}(R, Y | \Theta)}{\sum_{\tilde{Y}} \mathbb{P}(R, \tilde{Y} | \Theta)} \\ &= \underbrace{\frac{\mathbb{P}(R, Y | \Theta)}{\sum_{\tilde{Y}: c(\tilde{Y}) = c(Y)} \mathbb{P}(R, \tilde{Y} | \Theta)}}_{\mathbb{P}(Y | c(Y), R, \Theta)} \cdot \underbrace{\frac{\sum_{\tilde{Y}: c(\tilde{Y}) = c(Y)} \mathbb{P}(R, \tilde{Y} | \Theta)}{\sum_{c=0}^C \sum_{\tilde{Y}: c(\tilde{Y}) = c} \mathbb{P}(R, \tilde{Y} | \Theta)}}_{\mathbb{P}(c(Y) | R, \Theta)}. \end{aligned}$$

By the structure of our prior it follows that

$$(4) \quad \begin{aligned} \mathbb{P}(Y | c(Y), R, \Theta) &= \frac{\mathbb{P}(R | Y, \Theta) \mathbb{P}(Y)}{\sum_{\tilde{Y}: c(\tilde{Y}) = c(Y)} \mathbb{P}(R | \tilde{Y}, \Theta) \mathbb{P}(\tilde{Y})} \\ &= \frac{\mathbb{P}(R | Y, \Theta)}{\sum_{\tilde{Y}: c(\tilde{Y}) = c(Y)} \mathbb{P}(R | \tilde{Y}, \Theta)}, \end{aligned}$$

and

$$(5) \quad \begin{aligned} \mathbb{P}(c(Y) | R, \Theta) &= \frac{\sum_{\tilde{Y}: c(\tilde{Y}) = c(Y)} \mathbb{P}(R | \tilde{Y}, \Theta) \mathbb{P}(\tilde{Y})}{\sum_{c=0}^C \sum_{\tilde{Y}: c(\tilde{Y}) = c} \mathbb{P}(R | \tilde{Y}, \Theta) \mathbb{P}(\tilde{Y})} \\ &= \frac{\sum_{\tilde{Y}: c(\tilde{Y}) = c(Y)} \mathbb{P}(R | \tilde{Y}, \Theta) \mathbb{P}(\tilde{Y} | c(\tilde{Y})) \mathbb{P}(c(\tilde{Y}))}{\sum_{c=0}^C \sum_{\tilde{Y}: c(\tilde{Y}) = c} \mathbb{P}(R | \tilde{Y}, \Theta) \mathbb{P}(\tilde{Y} | c(\tilde{Y})) \mathbb{P}(c(\tilde{Y}))}. \end{aligned}$$

This decomposition suggests a good approach to sampling from $\mathbb{P}(Y | R, \Theta)$: we first sample $c(Y)$ according to $\mathbb{P}(c(Y) | R, \Theta)$ and then sample Y given the number of binding sites, according to $\mathbb{P}(Y | c(Y), R, \Theta)$.

As we will see next, we need to work more to obtain a centroid estimator for the binding sites: we need to establish a hierarchical inferential structure by first finding centroids for $c(Y) = 1, \dots, C$ and then proceed to estimate a global centroid. To this end we find $\mathbb{P}(c(Y) | R, \Theta)$ and then compute marginal posteriors $\mathbb{P}(Y_k | c(Y), R, \Theta)$.

3.1. Marginal posterior on $c(Y)$. From Equations 4 and 5 we observe that we need to compute $\sum_{\tilde{Y}: c(\tilde{Y}) = c} \mathbb{P}(R | \tilde{Y}, \Theta)$ up to a constant to find both conditional posteriors of $c(Y)$

and Y and thus the posterior $\mathbb{P}(Y | R, \Theta)$. Let us now denote by $R_{i:j}$ the subsequence of R from positions i to j and by $Y_{i:j}$ the binding sites in Y between i and j —that is, all Y_k such that $i \leq Y_k \leq j - L + 1$. If we then define *forward sums*

$$(6) \quad F_{c,j} \doteq \frac{\sum_{\tilde{Y}_{1:j}:c(\tilde{Y}_{1:j})=c} \mathbb{P}(R_{1:j} | \tilde{Y}_{1:j}, \Theta)}{\prod_{i=1}^j \prod_{s \in \mathcal{S}} \theta_{0,s}^{I(R_i=s)}}$$

we have that $\sum_{\tilde{Y}:c(\tilde{Y})=c} \mathbb{P}(R | \tilde{Y}, \Theta) \propto F_{c,n}$. To further simplify the notation, let us define

$$\lambda(j; \Theta) = \prod_{i=1}^L \prod_{s \in \mathcal{S}} \left(\frac{\theta_{i,s}}{\theta_{0,s}} \right)^{I(R_{j-1+i}=s)},$$

the composition ratio between motif and background for a binding site starting at j .

The forward sums $F_{c,j}$ can be computed recursively,

$$(7) \quad F_{c,j} = F_{c,j-1} + F_{c-1,j-L} \lambda(j-L+1; \Theta),$$

by considering two options for the tail of the sequence: either having a background position—and hence the first summand above—or by having a binding site on the last L positions—and thus requiring the second summand.

Thus, we have

$$(8) \quad \mathbb{P}(c(Y) | R, \Theta) = \frac{F_{c(Y),n} \binom{n-c(Y)(L-1)}{c(Y)}^{-1} \mathbb{P}(c(Y))}{\sum_{c=0}^C F_{c,n} \binom{n-c(L-1)}{c}^{-1} \mathbb{P}(c(Y) = c)},$$

which yields a straightforward way to sample the posterior $c(Y)$ conditional on Θ .

3.2. Marginal posterior on Y_k given $c(Y)$. To compute $\mathbb{P}(Y_k | c(Y), R, \Theta)$ we now need backward sums. We can define them analogously to the forward sums:

$$(9) \quad B_{c,j} \doteq \frac{\sum_{\tilde{Y}_{j:n}:c(\tilde{Y}_{j:n})=c} \mathbb{P}(R_{j:n} | \tilde{Y}_{j:n}, \Theta)}{\prod_{i=j}^n \prod_{s \in \mathcal{S}} \theta_{0,s}^{I(R_i=s)}},$$

and hence $\sum_{\tilde{Y}:c(\tilde{Y})=c} \mathbb{P}(R | \tilde{Y}, \Theta) \propto B_{c,1}$, as expected. Moreover, by a similar argument to the previous subsection, we also have that the backward sums are recursive:

$$(10) \quad B_{c,j} = B_{c,j+1} + B_{c-1,j+L} \lambda(j; \Theta).$$

Having forward and backward sums enable us to readily compute the marginal posterior on Y_k conditional on $c(Y)$: since

$$\begin{aligned} \mathbb{P}(Y_k | c(Y) = c, R, \Theta) &= \sum_{Y_1, \dots, Y_{k-1}, Y_{k+1}, \dots, Y_c} \mathbb{P}(Y | c(Y) = c, R, \Theta) \\ &= \sum_{Y_1, \dots, Y_{k-1}, Y_{k+1}, \dots, Y_c} \frac{\mathbb{P}(R | Y, \Theta)}{\sum_{\tilde{Y}:c(\tilde{Y})=c} \mathbb{P}(R | \tilde{Y}, \Theta)}, \end{aligned}$$

and

$$\sum_{Y_1, \dots, Y_{k-1}, Y_{k+1}, \dots, Y_c} \mathbb{P}(R | Y, \Theta) = \sum_{Y_1, \dots, Y_{k-1}} \mathbb{P}(R_{1:Y_{k-1}} | Y_{1:Y_{k-1}}, \Theta) \cdot \mathbb{P}(R_{Y_k:Y_{k+L}-1} | Y_{Y_k:Y_{k+L}-1}, \Theta) \cdot \sum_{Y_{k+1}, \dots, Y_c} \mathbb{P}(R_{Y_k+L:n} | Y_{Y_k+L:n}, \Theta),$$

and thus

$$(11) \quad \mathbb{P}(Y_k | c(Y) = c, R, \Theta) = \frac{F_{k-1, Y_{k-1}} \lambda(Y_k; \Theta) B_{c-k, Y_k+L}}{\sum_{\tilde{Y}_k=(k-1)L}^{n-(c-k+1)L+1} F_{k-1, \tilde{Y}_k-1} \lambda(\tilde{Y}_k; \Theta) B_{c-k, \tilde{Y}_k+L}}.$$

Note that

$$\begin{aligned} \frac{\sum_{\tilde{Y}:c(\tilde{Y})=c} \mathbb{P}(R | \tilde{Y}, \Theta)}{\prod_{i=1}^n \prod_{s \in \mathcal{S}} \theta_{0,s}^{I(R_i=s)}} &= F_{c,n} = B_{c,1} \\ &= \sum_{\tilde{Y}_k=(k-1)L}^{n-(c-k+1)L+1} F_{k-1, \tilde{Y}_k-1} \lambda(\tilde{Y}_k; \Theta) B_{c-k, \tilde{Y}_k+L}, \end{aligned}$$

for $k = 1, \dots, c$.

Before discussing posterior inference we summarize the results of this section in Algorithm 1.

Algorithm 1 Computes $\mathbb{P}(c(Y) | R, \Theta)$ and $\mathbb{P}(Y_k | c(Y), R, \Theta)$ for $k = 1, \dots, c(Y)$.

Step 1. (*Initialize*) Set $F_{0,0} = B_{0,n+1} = F_{0,j} = B_{0,j} = 1$ for $j = 1, \dots, n$; for $c = 1, \dots, C$, set $F_{c,j} = 0$ when $j < cL$ and $B_{c,j} = 0$ when $j > n - cL + 1$.

Step 2. (*Compute forward sums*) For $c = 1, \dots, C$ and $j = cL + 1, \dots, n$ do: set $F_{c,j}$ as in Equation 7,

$$F_{c,j} = F_{c,j-1} + F_{c-1,j-L} \lambda(j-L+1; \Theta)$$

Step 3. (*Compute $\mathbb{P}(c(Y) | R, \Theta)$*) For $c = 0, \dots, C$ do: compute marginal posterior $c(Y)$ as in Equation 8,

$$\mathbb{P}(c(Y) = c | R, \Theta) = \frac{F_{c,n} \binom{n-c(L-1)}{c}^{-1} \mathbb{P}(c(Y) = c)}{\sum_{\tilde{c}=0}^C F_{\tilde{c},n} \binom{n-\tilde{c}(L-1)}{\tilde{c}}^{-1} \mathbb{P}(c(Y) = \tilde{c})}$$

Step 4. (*Compute backward sums*) For $c = 1, \dots, C$ and $j = n - cL, \dots, 1$ do: set $B_{c,j}$ as in Equation 10,

$$B_{c,j} = B_{c,j+1} + B_{c-1,j+L} \lambda(j; \Theta)$$

Step 5. (*Compute $\mathbb{P}(Y_k | c(Y), R, \Theta)$*) For $c = 1, \dots, C$, $k = 1, \dots, c$, and $Y_k = (k-1)L + 1, \dots, n - (c-k+1)L + 1$ do: compute marginal posterior Y_k given $c(Y)$ as in Equation 11,

$$\mathbb{P}(Y_k | c(Y) = c, R, \Theta) = F_{k-1, Y_{k-1}} \lambda(Y_k; \Theta) B_{c-k, Y_k+L} / F_{c,n}$$

3.3. Posterior Inference. In contrast to the one binding site case from last section, posterior inference is more difficult since comparing configurations with different number of binding sites is not amenable to a systematic approach. Our first approximation is to consider local estimators for each group of configurations with a fixed number of binding sites and then appeal to a triangle inequality:

$$H(Y, \hat{Y}) \leq H(Y, \hat{Y}_c) + H(\hat{Y}_c, \hat{Y}),$$

where Y is a configuration with c binding sites, \hat{Y}_c is the constrained estimator for all configurations with c binding sites, and \hat{Y} is the (overall) centroid estimator. Recall that for the centroid estimator we wish to find \tilde{Y} that minimizes

$$\mathbb{E}_{Y|R,\Theta}[H(\tilde{Y}, Y)] = \sum_{c=0}^C \sum_{Y:c(Y)=c} H(\tilde{Y}, Y) \mathbb{P}(Y | R, \Theta).$$

Using the triangle inequality for each group we then have

$$(12) \quad \mathbb{E}_{Y|R,\Theta}[H(\tilde{Y}, Y)] \leq \sum_{c=0}^C \sum_{Y:c(Y)=c} [H(\tilde{Y}, \tilde{Y}_c) + H(\tilde{Y}_c, Y)] \mathbb{P}(Y | R, \Theta) \\ = \sum_{c=0}^C \left[H(\tilde{Y}, \tilde{Y}_c) + \sum_{Y:c(Y)=c} H(\tilde{Y}_c, Y) \mathbb{P}(Y | c(Y) = c, R, \Theta) \right] \mathbb{P}(c(Y) = c | R, \Theta),$$

where \tilde{Y}_c is an arbitrary point in $\{Y : c(Y) = c\}$. Our task is now to find an estimator—let us still call it centroid—that minimizes the right-hand bound in Equation 12 above. This goal suggests a two-step strategy:

1. For each number of binding sites, $c = 1, \dots, C$, find the *local* centroids

$$(13) \quad \hat{Y}_c = \arg \min_{\tilde{Y}:c(\tilde{Y})=c} \mathbb{E}_{Y|c(Y)=c,R,\Theta}[H(\tilde{Y}, Y)]$$

as the \tilde{Y}_c in Equation 12.

2. Find the *global* centroid given the local centroids $\{\hat{Y}_c\}_{c=1}^C$,

$$(14) \quad \hat{Y} = \arg \min_{\tilde{Y}} \mathbb{E}_{c(Y)|R,\Theta}[H(\hat{Y}_{c(Y)}, \tilde{Y})].$$

We note that this strategy does not guarantee that the bound is minimized; the main goal here is computational convenience. Let us tackle each step of this heuristic next.

3.3.1. Local centroids. Even when the number of binding sites is fixed, minimizing the conditional posterior expectation of $H(\tilde{Y}, Y)$ can be challenging: we would still have to consider for each candidate configuration \tilde{Y} the posterior probability of configurations with all binding sites to the left of the first binding site in \tilde{Y} , in-between binding sites in \tilde{Y} , and so on. We adopt another approximation and decide to minimize a *paired* Hamming loss H_A where binding site positions are matched according to their order:

$$H_A(\tilde{Y}, Y) = \sum_{k=1}^{c(Y)} H_1(\tilde{Y}_k, Y_k),$$

where $H_1(\tilde{Y}_k, Y_k)$ is Hamming loss when comparing sequences with only one binding site at \tilde{Y}_k and Y_k , respectively, that is, $H_1(\tilde{Y}_k, Y_k) = 2 \max\{|\tilde{Y}_k - Y_k|, L\}$. From the definition we have that H_A upper bounds H : $H_A(\tilde{Y}, Y) \geq H(\tilde{Y}, Y)$. As a bad approximation example, if $\tilde{Y}_k = Y_{k+1}$ for $k = 1, \dots, c(Y) - 1$ then $H_A(\tilde{Y}, Y) = c(Y)L$, since each pair of binding sites \tilde{Y}_k and Y_k does not overlap, while $H(\tilde{Y}, Y) = 2L$ since only Y_1 and $\tilde{Y}_{c(Y)}$ are in disagreement with background.

The next result adapts Theorem 1 to yield the paired local centroids.

LEMMA 2. *If $\mathbb{P}_k(\cdot | c(Y) = c, R, \Theta)$ is the marginal conditional posterior on Y_k then the paired local centroids are*

$$\hat{Y}_c = \arg \max_{\tilde{Y}: c(\tilde{Y})=c} \sum_{k=1}^c G(\tilde{Y}_k, \cdot) * \mathbb{P}_k(\cdot | c(Y) = c, R, \Theta)$$

PROOF. In the same spirit of Theorem 1, we use the conditional estimator in Equation 13 with the paired loss H_A :

$$\begin{aligned} \hat{Y}_c &= \arg \min_{\tilde{Y}: c(\tilde{Y})=c} \mathbb{E}_{Y | c(Y)=c, R, \Theta} [H_A(\tilde{Y}, Y)] \\ &= \arg \min_{\tilde{Y}: c(\tilde{Y})=c} \sum_{Y: c(Y)=c} \sum_{k=1}^c H_1(\tilde{Y}_k, Y_k) \mathbb{P}(Y | c(Y) = c, R, \Theta) \\ &= \arg \min_{\tilde{Y}: c(\tilde{Y})=c} \sum_{k=1}^c \sum_{Y_k=(k-1)L+1}^{n-(c-k+1)L+1} H_1(\tilde{Y}_k, Y_k) \mathbb{P}(Y_k | c(Y) = c, R, \Theta) \\ &= \arg \max_{\tilde{Y}: c(\tilde{Y})=c} \sum_{k=1}^c \sum_{Y_k=\max\{(k-1)L+1, \tilde{Y}_k-L\}}^{\min\{n-(c-k+1)L+1, \tilde{Y}_k+L\}} G(\tilde{Y}_k, Y_k) \mathbb{P}(Y_k | c(Y) = c, R, \Theta) \\ &= \arg \max_{\tilde{Y}: c(\tilde{Y})=c} \sum_{k=1}^c G(\tilde{Y}_k, \cdot) * \mathbb{P}_k(\cdot | c(Y) = c, R, \Theta), \end{aligned}$$

and the result follows. \square

We can spot in Lemma 2 the familiar convolutions, but now with the marginal posteriors $\mathbb{P}(Y_k | c(Y), R, \Theta)$ and in a more restricted range. We have a nice characterization, but we still have to optimize a sum to obtain the local centroids; to this end we explore the same recursive structure that allowed us to compute forward and backward sums. Let us define $f(\tilde{Y}_k) \doteq G(\tilde{Y}_k, \cdot) * \mathbb{P}_k(\cdot | c(Y) = c, R, \Theta)$ as the convolution against the marginal posterior on Y_k ; then we should have

$$(15) \quad \max_{\tilde{Y}: c(\tilde{Y})=c} \sum_{k=1}^c f(\tilde{Y}_k) = \max_{\tilde{Y}_c=(c-1)L+1, \dots, n-cL+1} \left[f(\tilde{Y}_c) + \max_{\tilde{Y}_1, \dots, \tilde{Y}_{c-1}} \sum_{k=1}^{c-1} f(\tilde{Y}_k) \right].$$

This important observation allows us to obtain \hat{Y}_c using the dynamic programming approach listed in Algorithm 2, as Theorem 3 formalizes.

THEOREM 3. *Algorithm 2 correctly identifies the paired local centroids*

$$\hat{Y}_c = \arg \min_{\tilde{Y}: c(\tilde{Y})=c} \mathbb{E}_{Y | c(Y)=c, R, \Theta} [H_A(\tilde{Y}, Y)].$$

PROOF. From Lemma 2 we know that \hat{Y}_c is the argument of $\max_{\tilde{Y}: c(\tilde{Y})=c} \sum_{k=1}^c f(\tilde{Y}_k)$. The key device in Algorithm 2 is to exploit the recursion in Equation 15 to define $m_1(\tilde{Y}_1) = f(\tilde{Y}_1)$ and

$$(16) \quad m_k(\tilde{Y}_k) = f(\tilde{Y}_k) + \max_{\tilde{Y}_{k-1}=(k-2)L+1, \dots, \tilde{Y}_k-L} m_{k-1}(\tilde{Y}_{k-1}),$$

Algorithm 2 Find \hat{Y}_c using dynamic programming.

Construct partial maxima and backtrack pointers:

Step 1. Set $m_1(\tilde{Y}_1) = f(\tilde{Y}_1)$ for $\tilde{Y}_1 = 1, \dots, n - cL + 1$.

Step 2. For $k = 2, \dots, c$ and $\tilde{Y}_k = (k-1)L + 1, \dots, n - (c-k+1)L + 1$ do: set backtrack pointers

$$A_{k-1}(\tilde{Y}_k) = \arg \max_{\tilde{Y}_{k-1}=(k-2)L+1, \dots, \tilde{Y}_k-L} m_{k-1}(\tilde{Y}_{k-1}).$$

and set partial sum maximum m_k as

$$m_k(\tilde{Y}_k) = f(\tilde{Y}_k) + m_{k-1}(A_{k-1}(\tilde{Y}_k)).$$

Reconstruct centroid \hat{Y}_c using backtrack pointers:

Step 3. Set last binding site position:

$$\hat{Y}_{c,c} = \arg \max_{\tilde{Y}_c=(c-1)L+1, \dots, n-L+1} m_c(\tilde{Y}_c).$$

Note that, by construction, $\max_{\tilde{Y}:c(\tilde{Y})=c} \sum_{k=1}^c f(\tilde{Y}_k) = m_c(\hat{Y}_{c,c})$.

Step 4. For $k = c, \dots, 2$ do: recover the remainder of \hat{Y}_c by setting $\hat{Y}_{c,k-1} = A_{k-1}(\hat{Y}_{c,k})$.

for $k > 1$, to store partial sum maxima. Now it follows that

$$\max_{\tilde{Y}:c(\tilde{Y})=c} \sum_{k=1}^c f(\tilde{Y}_k) = \max_{\tilde{Y}_c=(c-1)L+1, \dots, n-cL+1} m_c(\tilde{Y}_c),$$

and so Step 3 must be correct. The correctness of Step 4 relies on the right specification of m in Steps 1 and 2; but these steps are a straightforward application of Equation 15 using the definition of m_1 and a formulation of Equation 16 based on the backtrack pointers A , and so the algorithm is correct. \square

We note that the paired local centroids minimize an expected posterior upper bound H_A on the loss H , and so the actual local centroid might not be attained. We expect, however, that for common cases in which the motif coverage $c(Y)L$ is much smaller than n that the bound is tight since H_A approximates H well and thus the two local centroids often coincide.

3.3.2. Global centroid. While the local centroids already convey information about the distribution of posterior mass in the space of binding site configurations, the end goal of the analysis is a point estimate that is, in itself, a good representative of the space. Following the strategy we outlined in the beginning of this section, we can further summarize the information in the local centroids by identifying a configuration \hat{Y} that minimizes the expected conditional Hamming loss, as in Equation 14. This approach, however, entails the same difficulties as defining the centroid based on all points in the space, and it is thus not treatable by a systematic approach—we are now just restricting the configurations to the local centroids.

The global centroid can be defined by direct enumeration of all possible configurations while keeping the minimizer of the expected conditional posterior loss, but this “brute-force” approach considers an exponential number of solutions. A simple heuristic is to restrict the

global centroid to be one of the local centroids,

$$(17) \quad \hat{Y} = \arg \min_{\tilde{Y} \in \{\hat{Y}_c\}_{c=0}^C} \mathbb{E}_{c(Y) | R, \Theta} [H(\hat{Y}_{c(Y)}, \tilde{Y})].$$

Another alternative is to just take as global centroid the local centroid of the modal number of binding sites, $\hat{Y} = \hat{Y}_{c^*}$, where $c^* \doteq \arg \max_{c=0, \dots, C} \mathbb{P}(c(Y) = c | R, \Theta)$. From now on we adopt the global centroid in Equation 17 for simplicity and, again, computational expediency.

Before we continue to our next example, let us remark that a *constrained*, on the number of binding sites, global centroid might be more computationally feasible since we are restricting the space of available configurations. For instance, consider the 1-global centroid,

$$\hat{Y}_o \doteq \arg \min_{\tilde{Y}: c(\tilde{Y})=1} \mathbb{E}_{Y | R, \Theta} [H(\tilde{Y}, Y)].$$

As when defining local centroids, we can approximate \hat{Y}_o using a paired loss, and since

$$\begin{aligned} \mathbb{E}_{Y | R, \Theta} [H_A(\tilde{Y}, Y)] &= \sum_{c=0}^C \sum_{Y: c(Y)=c} \sum_{k=1}^c H_1(\tilde{Y}, Y_k) \mathbb{P}(Y | R, \Theta) \\ &= \sum_{i=1}^n \sum_{c=0}^C \sum_{Y: c(Y)=c} \sum_{k=1}^c H_1(\tilde{Y}, i) \mathbb{P}(Y_k = i | R, \Theta) \\ &= \sum_{i=1}^n H_1(\tilde{Y}, i) \sum_{c=0}^C \sum_{Y: c(Y)=c} \sum_{k=1}^c \mathbb{P}(Y_k = i | R, \Theta) \\ &= \sum_{i=1}^n H_1(\tilde{Y}, i) P_c(i | R, \Theta), \end{aligned}$$

where

$$(18) \quad P_c(i | R, \Theta) \doteq \sum_{c=1}^C \sum_{Y: c(Y)=c} \sum_{k=1}^c \mathbb{P}(Y_k = i | R, \Theta),$$

we have that

$$\hat{Y}_o = \arg \min_{\tilde{Y}: c(\tilde{Y})=1} \mathbb{E}_{Y | R, \Theta} [H_A(\tilde{Y}, Y)] = \arg \max_{\tilde{Y}: c(\tilde{Y})=1} G(\tilde{Y}, \cdot) * P_c(\cdot | R, \Theta).$$

It is important to note that while the restriction of one binding site might seem artificial, the derivation of \hat{Y}_o is helpful in recognizing sequence regions that are likely to host binding sites. In fact, since P_c captures the posterior probability of having a binding site starting at each position, and considering the overlap gain G , the convolution of G and P_c highlights positions that have higher posterior probability of being covered by a binding site.

EXAMPLE 2. We revisit the same sequence from Example 1, but now allow for at most $C = \lfloor n/L \rfloor = 33$ binding sites, and adopt the prior given in Equation 3 with $b = 3$ and thus $p = 1 - b/n = 0.985$. Using Algorithm 1 we are able to compute the conditional marginal posteriors $\mathbb{P}(c(Y) | R, \Theta)$ and $\mathbb{P}(Y_k | c(Y), R, \Theta)$ for $k = 1, \dots, c(Y)$. These posterior

distributions yield the local centroids—according to Algorithm 2—and the global centroid from Equation 17. In Table 2 we list the marginal posterior $\mathbb{P}(c(Y) = c | R, \Theta)$ up to the smallest c such that $\mathbb{P}(c(Y) \leq c | R, \Theta) > 0.95$, along with the local centroids; the global centroid \hat{Y}_C is highlighted. Interestingly, the global centroid coincides with the local centroid from the modal number of binding sites.

TABLE 2
Centroids and marginal posterior distribution of number of binding sites. The global centroid and the modal number of binding sites are highlighted in bold.

c	\hat{Y}_c	$\mathbb{P}(c(Y) = c R, \Theta)$	$\mathbb{P}(c(Y) \leq c R, \Theta)$
0	—	0.014	0.014
1	36	0.075	0.089
2	36, 147	0.181	0.270
3	13, 36, 147	0.254	0.524
4	13, 36, 63, 147	0.233	0.757
5	13, 36, 63, 147, 167	0.147	0.904
6	3, 29, 36, 63, 147, 167	0.067	0.971

In Figure 2 we display the posterior probabilities of binding site coverage P_c from Equation 18, along with the convolutions that are needed to define the 1-global centroid $\hat{Y}_o = 36$. As can be seen, position 36 has a lot of support, being present in all the local centroids listed in Table 2; in fact, the probability of a binding site starting at position 36 is greater than 50%.

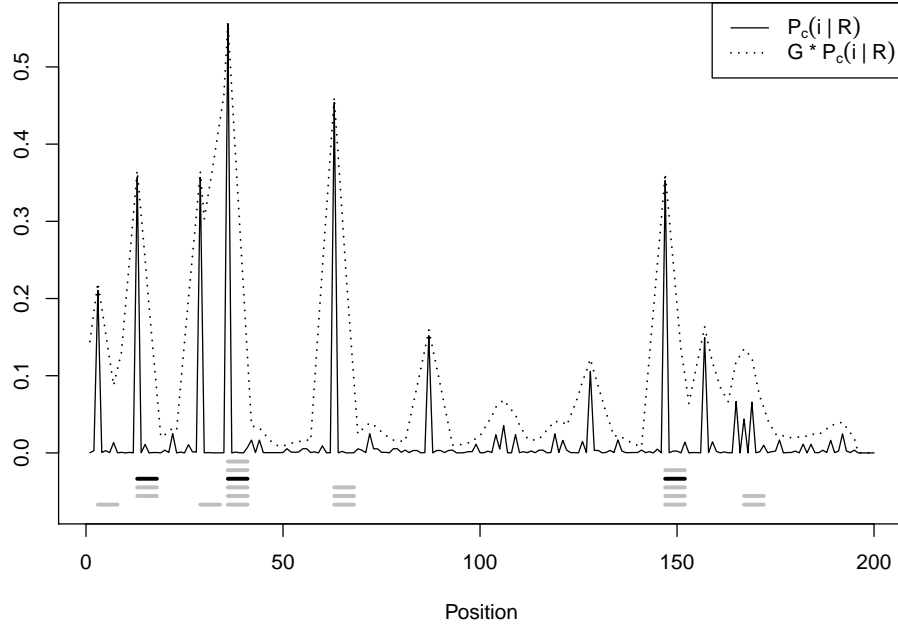


FIG 2. Posterior binding site coverage P_c in solid line and convolution $G * P_c$ in dotted line. Local centroids are listed below in gray; the global centroid is in black.

While P_c can provide us guidance for which positions are likely to start a binding site, using P_c to define local centroids can be misleading. For instance, we could expect that

the local centroid with three binding sites—the modal number of binding sites—would be, following a decreasing order on P_c , 36, 63, and 147. However, if we examine the marginal posteriors $\mathbb{P}(Y_k | c(Y) = 3, R, \Theta)$ in Figure 3 we realize that position 13 is favored over position 63 because, if $F_k \doteq G * \mathbb{P}_k(\cdot | c(Y) = 3, R, \Theta)$, $F_1(13) + F_2(36) > F_1(36) + F_2(63)$.

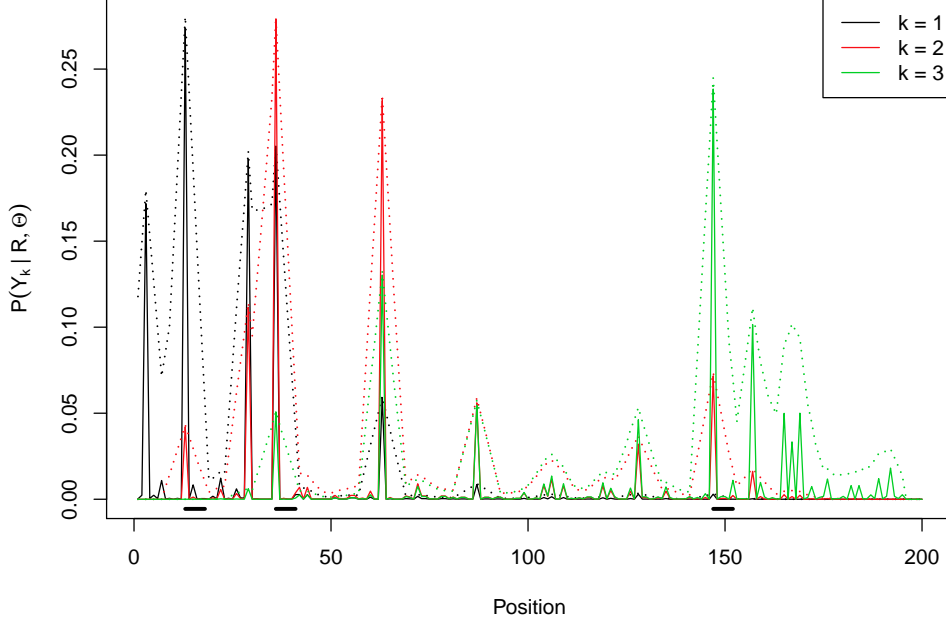


FIG 3. Marginal posterior distributions $\mathbb{P}(Y_k | c(Y) = 3, R, \Theta)$ in solid line and convolutions $G * \mathbb{P}(\cdot | c(Y), R, \Theta)$ in dotted line. The local centroid is displayed at the bottom.

4. Multiple sequences, multiple binding sites per sequence, random motif.

We are ready to address our model in broader generality: the dataset now comprises m sequences, $R = \{R_i\}_{i=1}^m$, and thus binding site configurations are also indexed by sequence, $Y = \{Y_i\}_{i=1}^m$. As before, we have that Y is independent of motif parameters Θ , but we further assume that sequences and configurations are conditionally independent given Θ :

$$(19) \quad \mathbb{P}(R, Y | \Theta) = \prod_{i=1}^m \mathbb{P}(R_i, Y_i | \Theta) = \prod_{i=1}^m \mathbb{P}(R_i | Y_i, \Theta) \mathbb{P}(Y_i).$$

Given Θ we would be able to apply the methods discussed this far to each sequence separately: compute forward and backward sums to obtain marginal posterior probabilities for each Y_i and then find local centroids and the i -th global centroid. We will, however, assume that Θ is random,

$$(20) \quad \theta_j \sim \text{Dir}(\alpha_j), \quad j = 0, 1, \dots, L,$$

independently, and we thus wish to also conduct inference on the background and motif compositions. This assumption, albeit more realistic, complicates matters, since the marginal

unconditioned posterior distributions of Y and Θ are not readily available; we are now required to estimate them before obtaining centroid estimates. To this end, we present next a Gibbs sampler (Geman and Geman, 1984; Liu, 2008) that draws Y_i for each sequence given Θ and then samples Θ conditional on the binding site configurations Y , similar to the approach in (Liu et al., 1995).

4.1. *Sampling Θ given Y and R .* Since the prior on Θ is conjugate, we should be able to sample Θ exactly from a Dirichlet distribution. From Equations 19 and 20 we have

$$\begin{aligned} \mathbb{P}(\theta_0 | Y, R) &\propto \left[\prod_{i=1}^m \prod_{s \in \mathcal{S}} \prod_{j \in BG_i} \theta_{0,s}^{I(R_{ij}=s)} \right] \left[\prod_{s \in \mathcal{S}} \theta_{0,s}^{\alpha_{0,s}-1} \right] \\ &= \prod_{s \in \mathcal{S}} \theta_{0,s}^{\sum_{i=1}^m \sum_{j \in BG_i} I(R_{ij}=s) + \alpha_{0,s} - 1}, \end{aligned}$$

and so $\theta_0 | Y, R \sim \text{Dir}(N_0(Y, R) + \alpha_0)$, where $N_0(Y, R) = \{N_{0,s}\}_{s \in \mathcal{S}}$ and

$$N_{0,s} = \sum_{i=1}^m \sum_{j \in BG_i} I(R_{ij} = s)$$

is the number of background positions across all sequences that have symbol s . Similarly, for the j -th position in the motif,

$$\begin{aligned} \mathbb{P}(\theta_j | Y, R) &\propto \left[\prod_{i=1}^m \prod_{s \in \mathcal{S}} \prod_{k=1}^{|Y_i|} \theta_{j,s}^{I(R_{i,Y_{ik}+j-1}=s)} \right] \left[\prod_{s \in \mathcal{S}} \theta_{j,s}^{\alpha_{j,s}-1} \right] \\ &= \prod_{s \in \mathcal{S}} \theta_{j,s}^{\sum_{i=1}^m \sum_{k=1}^{|Y_i|} I(R_{i,Y_{ik}+j-1}=s) + \alpha_{j,s} - 1}, \end{aligned}$$

and thus $\theta_j | Y, R \sim \text{Dir}(N_j(Y, R) + \alpha_j)$, with $N_j(Y, R) = \{N_{j,s}\}_{s \in \mathcal{S}}$ and

$$N_{j,s} = \sum_{i=1}^m \sum_{k=1}^{|Y_i|} I(R_{i,Y_{ik}+j-1} = s)$$

is the number of motif j -th positions across all sequences and binding sites that have symbol s .

4.2. *Sampling Y_i given Θ and R .* Each configuration Y_i for the i -th sequence is conditionally independent given Θ , so we can devise a sampling procedure that can be applied to each sequence in turn. To simplify the notation, let us drop the sequence index in what follows, that is, Y_i is Y , R_i is R , and so on. We will be following a similar approach to Sections 3.1 and 3.2, but instead of summing to obtain marginal distributions we will be sampling *exactly*.

To sample from the conditional posterior on Y , we first sample $c(Y) = c$ according to Equation 8 and then proceed to sample Y from its last, c -th binding site up to its first binding site. For this reason, this strategy is commonly referred to as “stochastic backtracking”, since it can be regarded as a stochastic version of Step 4 in Algorithms 1 and 2. Sampling Y is similar to the predictive update step in (Liu et al., 1995), which, on its turn, is based on a stochastic variation of expectation-maximization where missing data is imputed

(Tanner and Wong, 1987); however, here we exploit a hierarchical structure on $c(Y)$ and do not use the collapsing technique of Liu (1994).

Exploiting the conditional independence of the sequence configurations and Equation 4 the last binding site can be sampled using

$$\begin{aligned}
 \mathbb{P}(Y_c | c(Y), R, \Theta) &= \frac{\sum_{Y_1, \dots, Y_{c-1}} \mathbb{P}(R | Y, \Theta)}{\sum_{\tilde{Y}_c} \sum_{\tilde{Y}_1, \dots, \tilde{Y}_{c-1}} \mathbb{P}(R | \tilde{Y}, \Theta)} \\
 (21) \qquad &= \frac{F_{c-1, Y_{c-1}} \lambda(Y_c; \Theta)}{\sum_{\tilde{Y}_c=(c-1)L+1}^{n-L+1} F_{c-1, \tilde{Y}_{c-1}} \lambda(\tilde{Y}_c; \Theta)}.
 \end{aligned}$$

To sample the (intermediate) j -th binding site we use a similar expression:

$$\begin{aligned}
 \mathbb{P}(Y_j | Y_{j+1}, \dots, Y_c, c(Y), R, \Theta) &= \frac{\mathbb{P}(Y_j, \dots, Y_c, c(Y), R, \Theta)}{\sum_{\tilde{Y}_j} \mathbb{P}(\tilde{Y}_j, \dots, \tilde{Y}_c, c(Y), R, \Theta)} \\
 (22) \qquad &= \frac{\sum_{Y_1, \dots, Y_{j-1}} \mathbb{P}(R | Y, \Theta)}{\sum_{\tilde{Y}_j} \sum_{\tilde{Y}_1, \dots, \tilde{Y}_{j-1}} \mathbb{P}(R | \tilde{Y}, \Theta)} \\
 &= \frac{F_{j-1, Y_{j-1}} \lambda(Y_j; \Theta)}{\sum_{\tilde{Y}_j=(j-1)L+1}^{Y_{j+1}-L} F_{j-1, \tilde{Y}_{j-1}} \lambda(\tilde{Y}_j; \Theta)}.
 \end{aligned}$$

By making the convention that $Y_{c+1} = |R| + 1$ we can reduce Equation 21 to Equation 22. Moreover, note that Equation 22 implies that

$$\mathbb{P}(Y_j | Y_{j+1}, \dots, Y_c, c(Y), R, \Theta) = \mathbb{P}(Y_j | Y_{j+1}, c(Y), R, \Theta),$$

as expected.

We summarize the whole procedure in Algorithm 3. Note how Steps 1.1 to 1.3 are analogous to Steps 1 to 3 in Algorithm 1, and how Step 1.4 is a stochastic version of Step 4 in Algorithm 1: as previously stated, we are now sampling backwards instead of summing backwards. To obtain the centroids we follow the procedure described in Section 3.3, but adopting Monte Carlo estimates of the marginal posterior distributions, for $i = 1, \dots, m$,

$$\begin{aligned}
 \hat{\mathbb{P}}(c(Y_i) = c | R) &\approx \frac{1}{T} \sum_{t=1}^T I(c(Y_i^{(t)}) = c), \\
 \hat{\mathbb{P}}(Y_{ik} = j | c(Y_i) = c, R) &\approx \frac{\sum_{t=1}^T I(Y_{ik}^{(t)} = j) I(c(Y_i^{(t)}) = c)}{\sum_{t=1}^T I(c(Y_i^{(t)}) = c)}, \quad k = 1, \dots, c,
 \end{aligned}$$

where T is the number of samples.

EXAMPLE 3. For the random motif version of Example 2 we simulate $m = 20$ sequences of same length $n = 200$ using Θ from Table 1 and the prior for Y_i , $i = 1, \dots, m$, from Equation 3 with $p = 1 - 1/n = 0.995$.

We continue focusing on the inference of binding site configurations in the same sequence from previous examples, which is the first sequence in the simulated dataset. We assume a non-informative prior on Θ by setting $\alpha_{j,s} = 1$ for $s \in \mathcal{S}$ and $j = 0, \dots, L$; the prior on each sequence Y_i is the same prior from Example 2 with $p = 0.985$. Algorithm 3 is run for 10,000 iterations to guarantee convergence (diagnostics not shown.)

Algorithm 3 Gibbs sampler for $\mathbb{P}(Y, \Theta \mid R)$.

Set $\Theta^{(0)}$ arbitrarily. For $t = 1, \dots$ (until convergence) do:

Step 1. (*Sample $Y \mid \Theta, R$*) For each sequence $i = 1, \dots, m$, do: let $n = |R_i|$, $C = \lfloor n/L \rfloor$ and sample $Y_i \mid R_i, \Theta$.

Step 1.1. (*Initialize*) Set $F_{0,j} = 1$ for $j = 0, 1, \dots, n$ and for $c = 1, \dots, C$ set $F_{c,j} = 0$ when $j < cL$.

Step 1.2. (*Compute forward sums*) For $c = 1, \dots, C$ and $j = cL + 1, \dots, n$ do: set $F_{c,j}$ as in Equation 7,

$$F_{c,j} = F_{c,j-1} + F_{c-1,j-L} \lambda_i(j - L + 1; \Theta^{(t-1)}),$$

where λ_i uses R_i .

Step 1.3. (*Sample $c(Y_i^{(t)}) \mid R_i, \Theta^{(t-1)}$*) For $c = 0, \dots, C$ do: compute marginal posterior $c(Y_i)$ as in Equation 8 when applied to the i -th sequence,

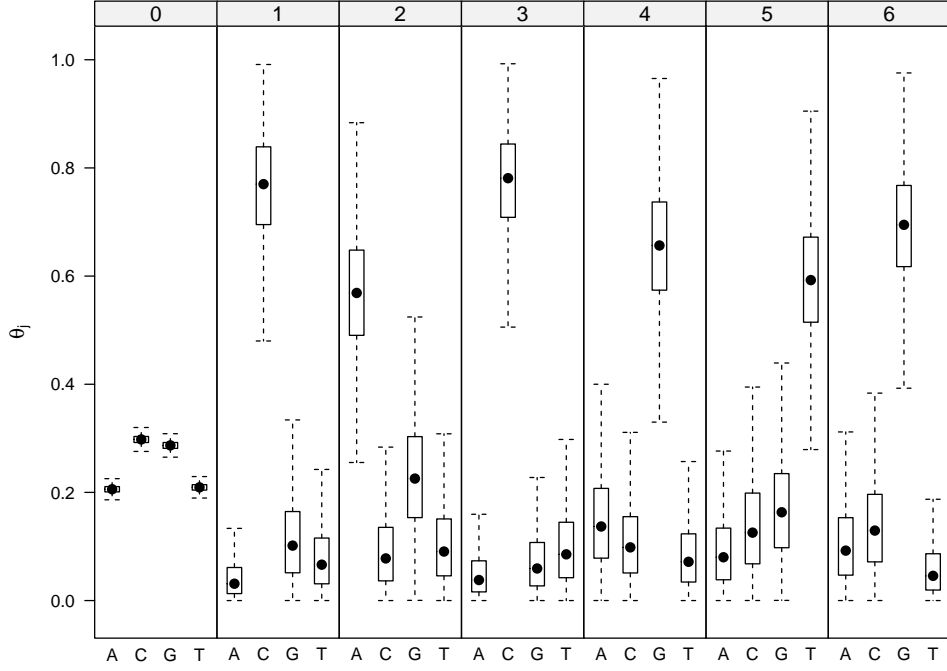
$$\mathbb{P}(c(Y_i) = c \mid R_i, \Theta^{(t-1)}) = \frac{F_{c,n} \binom{n-c(L-1)}{c}^{-1} \mathbb{P}(c(Y_i) = c)}{\sum_{\tilde{c}=0}^C F_{\tilde{c},n} \binom{n-\tilde{c}(L-1)}{\tilde{c}}^{-1} \mathbb{P}(c(Y_i) = \tilde{c})}$$

and sample $c^{(t)} \doteq c(Y_i^{(t)})$ according to $\mathbb{P}(c(Y_i) = c \mid R_i, \Theta^{(t-1)})$.

Step 1.4. (*Sample $Y_{ik}^{(t)} \mid c(Y_i^{(t)}) = c^{(t)}, R_i, \Theta^{(t-1)}$*) For $k = c^{(t)}, \dots, 1$ do: sample $Y_{ik}^{(t)}$ proportional to $F_{k-1, Y_{ik}^{(t)}-1} \lambda_i(Y_{ik}^{(t)}; \Theta^{(t-1)})$ as in Equation 22,

$$\begin{aligned} \mathbb{P}(Y_{ik}^{(t)} \mid Y_{i,k+1}^{(t)}, c(Y_i^{(t)}) = c^{(t)}, R_i, \Theta^{(t-1)}) &= \\ &= \frac{F_{k-1, Y_{ik}^{(t)}-1} \lambda_i(Y_{ik}^{(t)}; \Theta^{(t-1)})}{\sum_{\tilde{Y}_k = (k-1)L+1}^{Y_{i,k+1}^{(t)}-L} F_{k-1, \tilde{Y}_k-1} \lambda_i(\tilde{Y}_k; \Theta^{(t-1)})} \end{aligned}$$

Step 2. (*Sample $\Theta \mid Y, R$*) For $j = 0, \dots, L$ compute $N_j(Y^{(t)}, R)$ and then sample $\theta_j^{(t)} \mid Y^{(t)}, R \sim \text{Dir}(N_j(Y^{(t)}, R) + \alpha_j)$.

FIG 4. Boxplots of MCMC samples for Θ (outliers are not shown.)

The marginal posterior distribution of Θ can be assessed in Figure 4. Since most positions in the sequences are background sequences θ_0 has very small posterior variances. Also note that the canonical palindromic E-box motif, with consensus **CACGTG**, is recovered.

The procedure is now similar to what we presented in Example 2; the main difference is that the marginal posterior distributions are estimated from the MCMC samples. Table 3 lists the estimated marginal posterior distribution of the number of binding sites, the local and global centroids. The global centroid does not coincide with the local centroid for the modal number of binding sites. Moreover, the local centroids here are different from the (conditional) local centroids in Example 2, most likely due to the randomness of Θ being taken into account.

TABLE 3

Centroids and estimated marginal posterior distribution of number of binding sites. The global centroid and the modal number of binding sites are highlighted in bold.

c	\hat{Y}_c	$\hat{\mathbb{P}}(c(Y) = c R, \Theta)$	$\hat{\mathbb{P}}(c(Y) \leq c R, \Theta)$
0	—	0.026	0.026
1	29	0.107	0.133
2	29, 167	0.210	0.343
3	29, 63, 167	0.274	0.617
4	13, 36, 147, 167	0.201	0.818
5	13, 29, 63, 147, 167	0.120	0.938
6	13, 29, 36, 63, 147, 167	0.046	0.984

Figure 5 displays the estimated P_c , $G * P_c$, and the centroids. We see that compared to Example 2 some posterior mass has shifted to positions 29 and to the group of positions 166, 167, and 168. Here we clearly see the advantage of a centroid estimator: $G * P_c$, and

later $G * \mathbb{P}_k(\cdot | R)$, gathers evidence of motif binding from nearby positions, yielding a better summary—according to our choice of loss function—of the distribution of posterior mass.

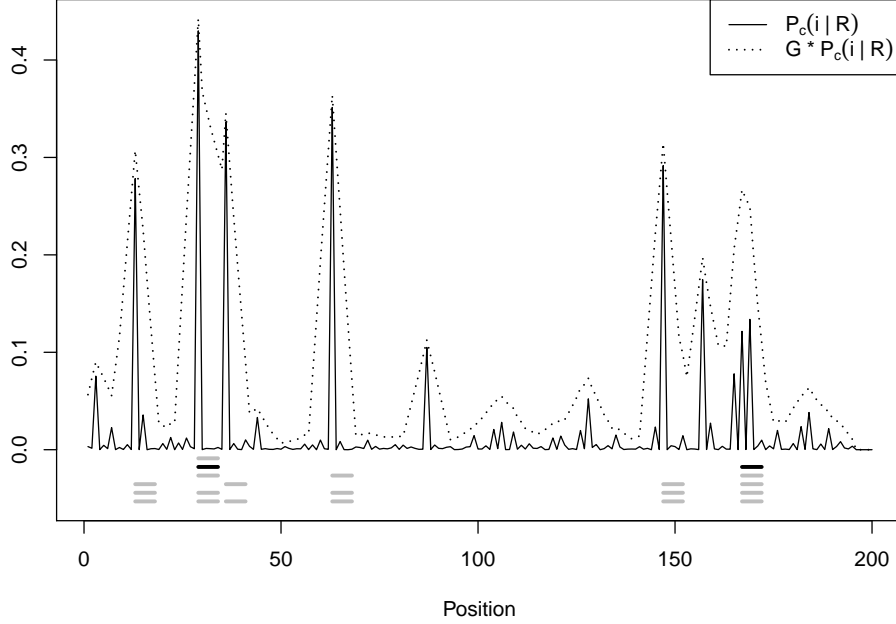


FIG 5. Estimated posterior binding site coverage P_c in solid line and convolution $G * P_c$ in dotted line. Local centroids are listed below in gray; the global centroid is in black.

The selection of position 167 in the second local centroid \hat{Y}_2 might seem puzzling since the peaks at positions 36, 63, and 147 hold higher coverage probabilities. Checking $\hat{\mathbb{P}}(Y_k | R)$ in Figure 6 helps dismiss any doubts: most of the support for these positions come from configurations with higher number of binding sites, as evidenced by the respective local centroids, but these configurations hold relatively low posterior mass. When $c(Y) = 2$, the prior on $Y_{2,2}$ assigns more posterior probability to higher positions, close to the end of the sequence, simply because there are more configurations for $Y_{2,2}$ on these positions. It is also important to notice that while none of the positions in the cluster 166–168 has higher marginal posterior mass than positions 63 and 147, the convolution $G * \hat{\mathbb{P}}_2(\cdot | R)$ is maximized at position 167, that is, the cluster when taken together has more support from the data, as weighted by G .

EXAMPLE 4. We end this section with an example from the real-world dataset in (Tompa et al., 2005), sequence set `yst02r`. The dataset contains $m = 4$ sequences each with $n = 500$ letters. We set $L = 16$ and adopt a non-informative prior on Θ , as in the previous example, and the prior on each Y_i , for the i -th sequence, from Equation 3 with $b = 3$ per thousand positions, so $p = 1 - 3/1000 = 0.997$. As in the previous example, 10,000 iterations suffice to reach convergence.

Let us focus on the second sequence. Figure 7 pictures the binding site coverage probabilities, along with the local centroids. The global centroid $\hat{Y}_C = \{85, 105, 169\}$ contains three binding sites, and it is also the local centroid for the modal number of binding sites, with

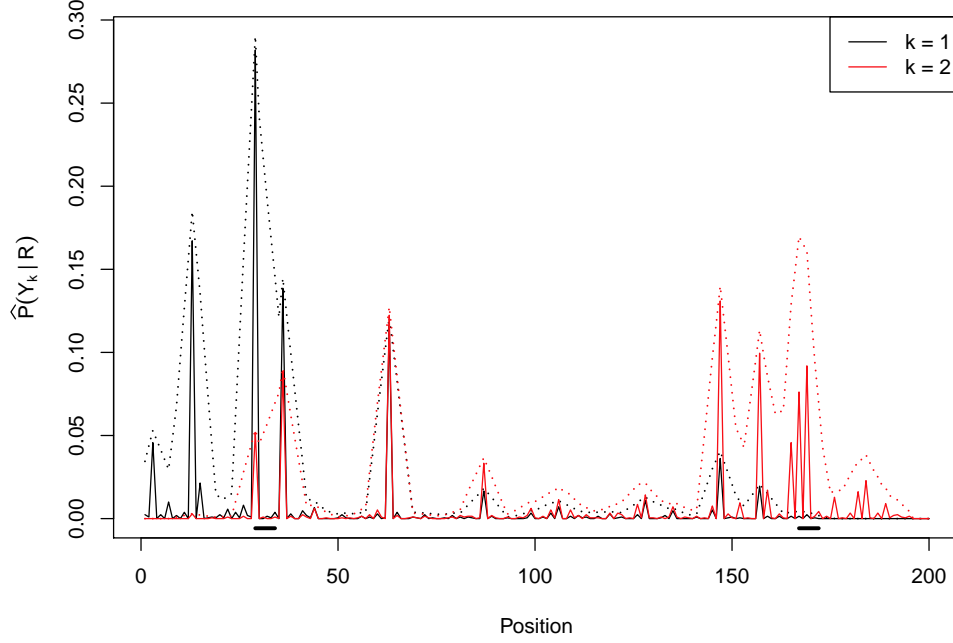


FIG 6. Estimated marginal posterior distributions $\mathbb{P}(Y_k | c(Y) = 2, R, \Theta)$ in solid line and convolutions $G * \mathbb{P}(\cdot | c(Y), R, \Theta)$ in dotted line. The local centroid is displayed at the bottom.

$\hat{\mathbb{P}}(c(Y) = 3 | R) = 0.32$. Since most of the posterior mass is concentrated in configurations with $c(Y) = 3$, the posterior profiles $\hat{\mathbb{P}}(Y_k | c(Y) = 3, R)$ are similar to P_c and are thus omitted.

From the MCMC samples we can produce the MAP estimate $\hat{Y}_M = \{86, 105, 174\}$ as the configuration with highest frequency among the samples: $\hat{\mathbb{P}}(\hat{Y}_M | R) = 0.032$. In fact, we can estimate the posterior probability of each sampled binding site configuration and then, using classic multidimensional scaling (Gower, 1966), visualize the estimated posterior distribution in Figure 8. It is interesting to note that the null configuration—that is, without binding sites—is also very likely with posterior probability 0.024. In contrast, the global centroid has very small posterior probability, close to 0.001; it sits, however, closer to configurations with high posterior mass, including the local centroids with one, two, and four binding sites.

To better assess how the centroid estimator is closer to a mean than a mode estimator, we plot the estimated posterior distribution of the generalized loss function H centered at both \hat{Y}_C and \hat{Y}_M in Figure 9. Since $\mathbb{E}_{Y|R}[H(\hat{Y}_M, Y)] = 42.40$ and $\mathbb{E}_{Y|R}[H(\hat{Y}_C, Y)] = 40.22$, we see that the binding sites in the centroid configuration are, on average, overlapping two extra positions with the binding sites in all the configurations when compared to the MAP estimate's binding sites. Both estimates are fairly similar, but the centroid reminds us that placing the third binding site at position 169, instead of 174, yields an unlikely configuration, but with a higher chance of overlapping with binding sites in positions 160–175 that have high posterior probability. In the context of Figures 8 and 9, the centroid places itself between two clusters that concentrate posterior mass: one with configurations Y such that $25 \leq H(\hat{Y}_C, Y) \leq 40$ and another with configurations further away, satisfying

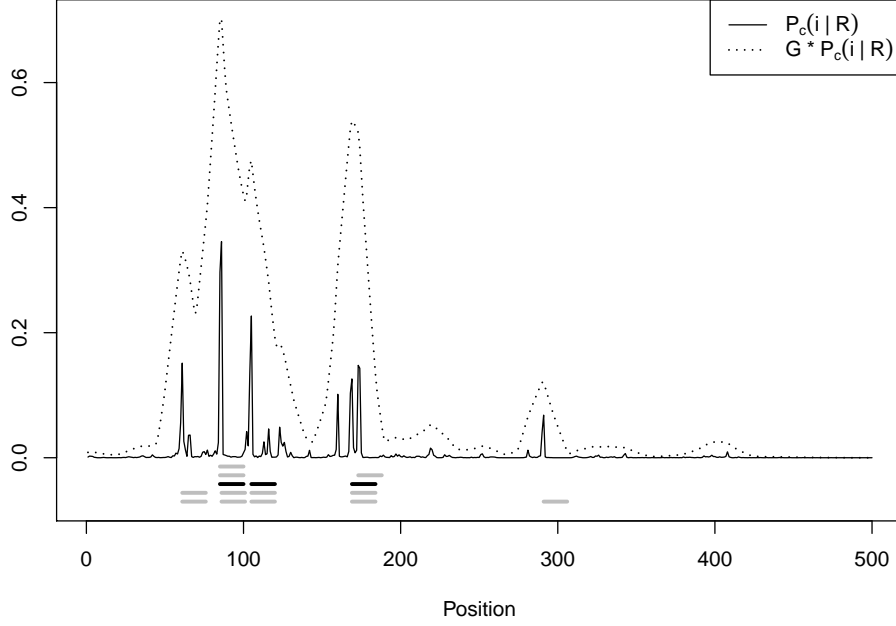


FIG 7. Estimated posterior binding site coverage P_c in solid line and convolution $G * P_c$ in dotted line for real-world dataset, second sequence. Local centroids are listed below in gray; the global centroid is in black.

$$40 \leq H(\hat{Y}_C, Y) \leq 50.$$

5. Discussion. In this paper we have presented a Bayesian approach, similar to the Gibbs motif sampler in (Lawrence et al., 1993; Liu et al., 1995), that jointly models motif and background compositions and binding site locations in a set of sequences. More importantly, we discuss and formalize an inferential procedure based on the centroid estimator proposed by Carvalho and Lawrence (2008). As in any Bayesian analysis, we wish to evaluate features of interest in a model based on their posterior distribution; however, if we are required to pick a representative configuration, a point in the parameter space, then a principled approach is to elect a loss function and conduct formal statistical decision analysis. In this sense, by exploring a more refined loss function that depends on position-wise comparisons between sequence states—background or motif positions—we are able to identify a better representative of the posterior space of binding site configurations. As pointed out in (Carvalho and Lawrence, 2008), the centroid estimator better accounts for the distribution of posterior mass; it is more similar to a median than to a mode, and can thus offer better predictive resolution than the MAP estimator (Barbieri and Berger, 2004). When applied to motif discovery, the centroid estimator captures information in the vicinity of binding site positions through a convolution in marginal posterior distributions of binding sites.

Given the combinatorial number of possible configurations in the parameter space it is not feasible to identify the centroid estimate through enumeration or even a systematic approach. Yet, we devise an approximative scheme that efficiently optimizes an upper bound on the posterior expected loss and thus provides a related centroid. Despite its heuristic nature, the proposed method has another advantage besides computational convenience: it

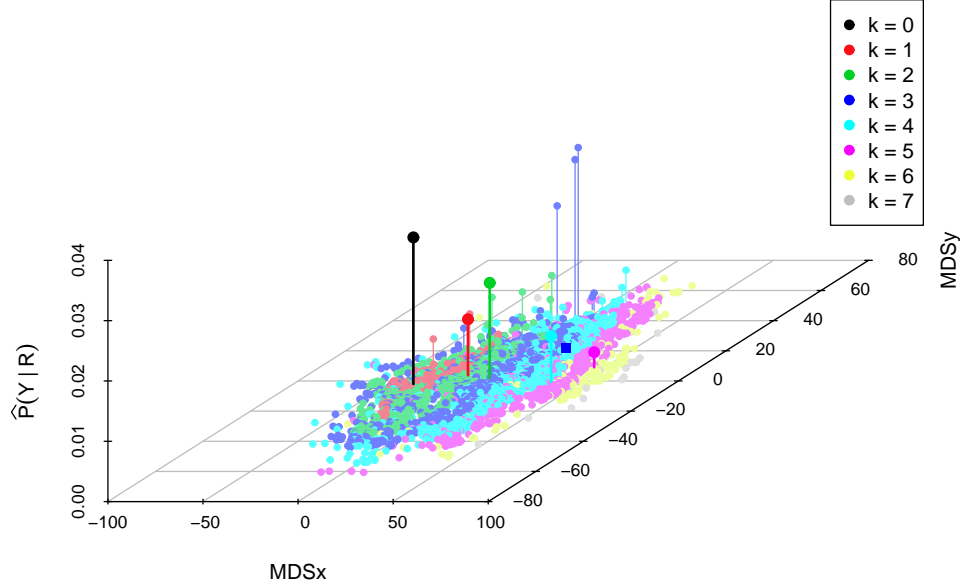


FIG 8. Estimated posterior distribution of configurations Y based on MCMC samples and projected using multidimensional scaling. The colors code configurations with different number of binding sites. Bold points mark local centroids, while a square (bold) point highlights the global centroid.

allows for an informative depiction of the posterior distribution on binding site configurations. First, when defining the local centroids, we are able to assess the contributions from each binding site through their marginal posterior distributions conditional on the number of binding sites, and, in particular, through the convolution of these marginal profiles with the gain filter; secondly, when finding the global centroid we explore the marginal posterior distribution on the number of binding sites. Moreover, other representations might be helpful in understanding the distribution of posterior mass, as in the use of P_c (in Equation 18) to pinpoint the 1-global centroid and measure the overall support of the configurations to a binding site at some specific position in the sequence. These comments are in the spirit of an estimator being also a communicator of the posterior space and the particular choice of prior distribution (see Berger, 1985, Section 4.10).

It is important to note that even when the model is accurate, a poor inference might fail in recovering relevant features of the space. In Example 2, the MAP estimate is the null configuration, while the centroid indicates three binding sites that represent a group of configurations that jointly pool significant posterior mass. It is also common that the posterior distribution is too complex to be reasonably captured by a single representative; in this case the expected posterior loss could also be used to partition the space and further define additional representatives as conditional estimates on each subspace. This is a direction of work that warrants interest and that we intend to follow next.

Further improvements can be obtained by specifying a more complex model that accounts, for example, for higher order Markov chains with more states for the background, as in (Roth et al., 1998; Liu et al., 2001), phylogenetic profiles (Newberg et al., 2007), structural

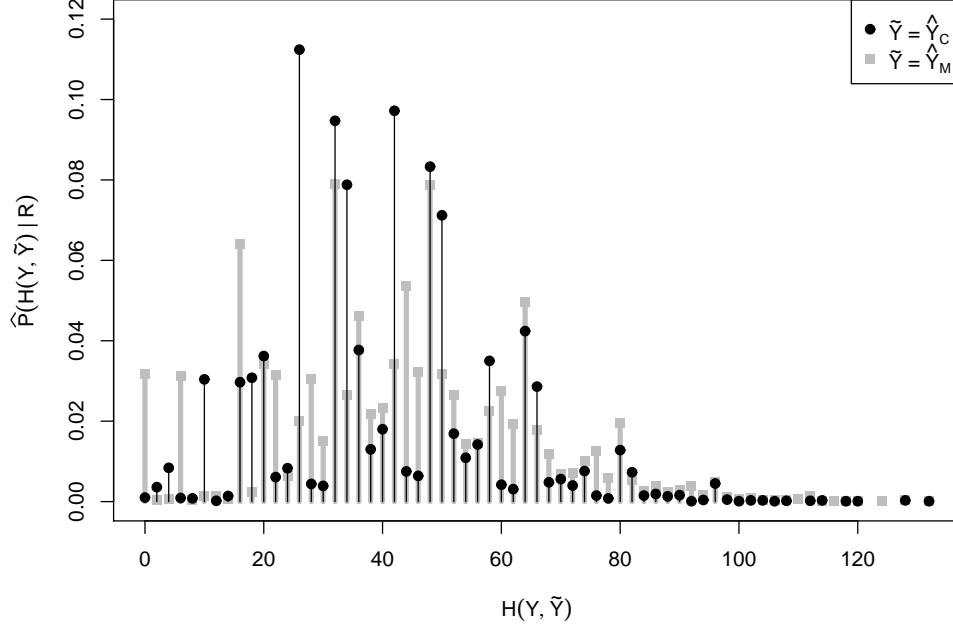


FIG 9. Estimated posterior distribution of loss function centered at \tilde{Y} for the MAP ($\tilde{Y} = \hat{Y}_M$) and centroid ($\tilde{Y} = \hat{Y}_C$) estimates.

information (Xing and Karp, 2004), a variable motif length, or dependency among motif positions. As pointed out by Hu et al. (2005), motif discovery using sequence only is well known for low signal-to-noise ratio; future extensions would also incorporate other data sources, such as gene expression or ChIP-Seq data, to increase the signal-to-noise ratio.

Acknowledgements. The author would like to thank Antonio Gomes for the helpful discussions and comments in the text.

References.

- Bailey, T. and C. Elkan (1995). Unsupervised learning of multiple motifs in biopolymers using expectation maximization. *Machine learning* 21(1), 51–80.
- Barbieri, M. and J. Berger (2004). Optimal predictive model selection. *The Annals of Statistics* 32(3), 870–897.
- Berger, J. (1985). *Statistical decision theory and Bayesian analysis*. Springer.
- Besag, J. (1986). On the statistical analysis of dirty pictures. *Journal of the Royal Statistical Society. Series B (Methodological)* 48(3), 259–302.
- Carvalho, L. and C. Lawrence (2008). Centroid estimation in discrete high-dimensional spaces with applications in biology. *Proceedings of the National Academy of Sciences of the United States of America* 105(9), 3209.
- Dempster, A., N. Laird, and D. Rubin (1977). Maximum likelihood from incomplete data via the EM algorithm. *Journal of the Royal Statistical Society. Series B (Methodological)* 39(1), 1–38.
- Ding, Y., C. Chan, and C. Lawrence (2005). RNA secondary structure prediction by centroids in a Boltzmann weighted ensemble. *RNA* 11(8), 1157–1166.
- Geman, S. and D. Geman (1984). Stochastic relaxation, Gibbs distributions, and the Bayesian restoration of images. *Pattern Analysis and Machine Intelligence, IEEE Transactions on* 6(6), 721–741.
- Gower, J. (1966). Some distance properties of latent root and vector methods used in multivariate analysis. *Biometrika* 53(3-4), 325–338.

- GuhaThakurta, D. (2006). Computational identification of transcriptional regulatory elements in DNA sequence. *Nucleic acids research* 34(12), 3585–3598.
- Hu, J., B. Li, and D. Kihara (2005). Limitations and potentials of current motif discovery algorithms. *Nucleic acids research* 33(15), 4899–4913.
- Lawrence, C., S. Altschul, M. Boguski, J. Liu, A. Neuwald, and J. Wootton (1993). Detecting subtle sequence signals: a Gibbs sampling strategy for multiple alignment. *Science* 262(5131), 208–214.
- Liu, J. (1994). The collapsed Gibbs sampler in Bayesian computations with applications to a gene regulation problem. *Journal of the American Statistical Association* 89(427), 958–966.
- Liu, J. (2008). *Monte Carlo strategies in scientific computing*. Springer Verlag.
- Liu, J., A. Neuwald, and C. Lawrence (1995). Bayesian models for multiple local sequence alignment and Gibbs sampling strategies. *Journal of the American Statistical Association* 90(432), 1156–1170.
- Liu, X., D. Brutlag, and J. Liu (2001). BioProspector: discovering conserved DNA motifs in upstream regulatory regions of co-expressed genes. In *Pac Symp Biocomput*, Volume 6, pp. 127–138.
- MacIsaac, K. and E. Fraenkel (2006). Practical strategies for discovering regulatory DNA sequence motifs. *PLoS computational biology* 2(4), e36.
- Murrea, C., P. S. McCawa, and D. Baltimorea (1989). A new DNA binding and dimerization motif in immunoglobulin enhancer binding, daughterless, MyoD, and Myc proteins. *Cell* 56(5), 777–783.
- Neuwald, A., J. Liu, and C. Lawrence (1995). Gibbs motif sampling: detection of bacterial outer membrane protein repeats. *Protein science* 4(8), 1618–1632.
- Newberg, L., W. Thompson, S. Conlan, T. Smith, L. McCue, and C. Lawrence (2007). A phylogenetic Gibbs sampler that yields centroid solutions for cis-regulatory site prediction. *Bioinformatics* 23(14), 1718–1727.
- Pavesi, G., P. Mereghetti, G. Mauri, and G. Pesole (2004). Weeder Web: discovery of transcription factor binding sites in a set of sequences from co-regulated genes. *Nucleic acids research* 32(suppl 2), W199–W203.
- Pevzner, P., S. Sze, et al. (2000). Combinatorial approaches to finding subtle signals in DNA sequences. In *Proceedings of the Eighth International Conference on Intelligent Systems for Molecular Biology*, Volume 8, pp. 269–278.
- Régner, M. and A. Denise (2004). Rare events and conditional events on random strings. *Discrete Mathematics and Theoretical Computer Science* 6(2), 191–214.
- Roth, F., J. Hughes, P. Estep, and G. Church (1998). Finding DNA regulatory motifs within unaligned noncoding sequences clustered by whole-genome mRNA quantitation. *Nature biotechnology* 16(10), 939–945.
- Sandve, G. and F. Drablos (2006). A survey of motif discovery methods in an integrated framework. *Biol Direct* 1(11).
- Stormo, G. (2000). DNA binding sites: representation and discovery. *Bioinformatics* 16(1), 16–23.
- Tanner, M. and W. Wong (1987). The calculation of posterior distributions by data augmentation. *Journal of the American statistical Association* 82(398), 528–540.
- Thijs, G., K. Marchal, M. Lescot, S. Rombauts, B. De Moor, P. Rouze, and Y. Moreau (2002). A Gibbs sampling method to detect overrepresented motifs in the upstream regions of coexpressed genes. *Journal of Computational Biology* 9(2), 447–464.
- Thompson, W., L. Newberg, S. Conlan, L. McCue, and C. Lawrence (2007). The Gibbs centroid sampler. *Nucleic acids research* 35(suppl 2), W232–W237.
- Tompa, M., N. Li, T. Bailey, G. Church, B. De Moor, E. Eskin, A. Favorov, M. Frith, Y. Fu, W. Kent, et al. (2005). Assessing computational tools for the discovery of transcription factor binding sites. *Nature biotechnology* 23(1), 137–144.
- Xing, E. and R. Karp (2004). MotifPrototyper: a bayesian profile model for motif families. *Proceedings of the National Academy of Sciences of the United States of America* 101(29), 10523.

DEPARTMENT OF MATHEMATICS AND STATISTICS
 BOSTON UNIVERSITY
 BOSTON, MASSACHUSETTS 02215
 E-MAIL: lecarval@math.bu.edu